
Near-Isometric Binary Hashing for Large-scale Datasets

Amirali Aghazadeh

AMIRALI@RICE.EDU

Department of Electrical and Computer Engineering, Rice University, Houston, TX, USA

Andrew Lan

MR.LAN@SPARFA.COM

Department of Electrical and Computer Engineering, Rice University, Houston, TX, USA

Anshumali Shrivastava

ANSHUMALI@RICE.EDU

Department of Computer Science, Rice University, Houston, TX, USA

Richard Baraniuk

RICHB@RICE.EDU

Department of Electrical and Computer Engineering, Rice University, Houston, TX, USA

Abstract

We develop a scalable algorithm to learn binary hash codes for indexing large-scale datasets. Near-isometric binary hashing (NIBH) is a data-dependent hashing scheme that quantizes the output of a learned low-dimensional embedding to obtain a binary hash code. In contrast to conventional hashing schemes, which typically rely on an ℓ_2 -norm (i.e., average distortion) minimization, NIBH is based on a ℓ_∞ -norm (i.e., worst-case distortion) minimization that provides several benefits, including superior distance, ranking, and near-neighbor preservation performance. We develop a practical and efficient algorithm for NIBH based on column generation that scales well to large datasets. A range of experimental evaluations demonstrate the superiority of NIBH over ten state-of-the-art binary hashing schemes.

In this paper, we are interested in learning *near-isometric* binary embeddings, i.e., hash functions that preserve the distances between data points up to a given distortion in the Hamming space. More rigorously, let \mathcal{Z} and \mathcal{Y} denote metric spaces with metrics $d_{\mathcal{Z}}$ and $d_{\mathcal{Y}}$, respectively. An embedding $f : \mathcal{Z} \rightarrow \mathcal{Y}$ is called near-isometric (Plan & Vershynin, 2014, Def. 1.1) if, for *every* pair of data points $\mathbf{z}_i, \mathbf{z}_j \in \mathcal{Z}$, we have

$$d_{\mathcal{Z}}(\mathbf{z}_i, \mathbf{z}_j) - \gamma \leq d_{\mathcal{Y}}(f(\mathbf{z}_i), f(\mathbf{z}_j)) \leq d_{\mathcal{Z}}(\mathbf{z}_i, \mathbf{z}_j) + \gamma,$$

where γ is called the isometry constant. In words, f is near-isometric if and only if the entries of the *pairwise-distortion vector* containing the distance distortion between every pair of data points $|d_{\mathcal{Z}}(\mathbf{z}_i, \mathbf{z}_j) - d_{\mathcal{Y}}(f(\mathbf{z}_i), f(\mathbf{z}_j))|, \forall i > j$ do not exceed the isometry constant γ . A near-isometric embedding is approximately *distance-preserving* in that the distance between any pairs of data points in the embedded space \mathcal{Y} is approximately equal to their distance in the ambient space \mathcal{Z} (Hegde et al., 2015; Shaw & Jebara, 2007; Weinberger & Saul, 2006).

1. Introduction

Hashing, one of the primitive operations in large-scale systems, seeks a low-dimensional binary embedding of a high-dimensional data set. Such a binary embedding can increase the computational efficiency of a variety of tasks, including searching, learning, near-neighbor retrieval, etc. One of the fundamental challenges in machine learning is the development of efficient hashing algorithms that embed data points into compact binary codes while preserving the geometry of the original dataset.

The simplest and most popular binary hashing scheme, *random projection*, simply projects the data into a lower-dimensional (lower-bit) random subspace and then quantizes to binary values. Random projections are known to be near-isometric with high probability, due to the celebrated Johnson-Lindenstrauss (JL) lemma (Andoni et al., 2014; Datar et al., 2004; Plan & Vershynin, 2014). Algorithms based on the JL lemma belong to the family of probabilistic dimensionality reduction techniques; a notable example is locality sensitive hashing (LSH) (Andoni et al., 2014; Datar et al., 2004). Unfortunately, theoretical results on LSH state that the number of bits required to guarantee an isometric embedding can be as large as the number of data points (Datar et al., 2004; Plan & Vershynin, 2014). Even

in practice, LSH’s requirement on the number of bits is impractically high for indexing many real-world, large-scale datasets (Lv et al., 2007).

Consequently, several *data-dependent* binary hashing algorithms have been developed that leverage the structure of the data to learn compact binary codes. These methods enable a significant reduction in the number of bits required to index large-scale datasets compared to LSH; see (Wang et al., 2014) for a survey. However, learning compact binary codes that preserve the local geometry of the data remains challenging.

These data-dependent hashing algorithms focus on the choice of the distortion measure. Typically, after finding the appropriate distortion measure, the hash functions are learned by minimizing the *average* distortion, i.e., the ℓ_2 -norm of the pairwise-distortion vector, which sums the distortion among all pairwise distances with equal weights. *Binary reconstructive embedding* (BRE) (Kulis & Darrell, 2009), for example, uses an optimization algorithm to directly minimize the average distortion in the embedded space. *Spectral Hashing* (SH) (Weiss et al., 2009), *Anchor Graph Hashing* (AGH) (Liu et al., 2011), *Multidimensional Spectral Hashing* (MDSH) (Weiss et al., 2012), and *Scalable Graph Hashing* (SGH) (Jiang & Li, 2015) define notions of *similarity* based on a function of ℓ_2 -distance between the data points $\|\mathbf{z}_i - \mathbf{z}_j\|_2$ and use spectral methods to learn hash functions that keep similar points sufficiently close. Some other hashing algorithms first project the data onto its principal components, e.g., *PCA Hashing* (PCAH) (Jolliffe, 2002), which embeds with minimal average distortion, and then learn a rotation matrix to minimize the quantization error (Gong & Lazebnik, 2011) or balance the variance across the components (*Isotropic Hashing* (Iso-Hash) (Kong & Li, 2012)).

While minimizing the average distortion seems natural, this approach can sacrifice the preservation of certain pairwise distances in favor of others. As we demonstrate below, this can lead to poor performance in certain applications, such as the preservation of nearest neighbors. In response, in this paper, we develop a new data-driven hashing algorithm that minimizes the *worst-case* distortion among the pairwise distances, i.e., the ℓ_∞ -norm of the pairwise-distortion vector.

Figure 1 illustrates the advantages of minimizing the ℓ_∞ -norm of the pairwise-distortion vector instead of its ℓ_2 -norm. Consider three clusters of points in a two-dimensional space. We compute the optimal one-dimensional embeddings of the data points by minimizing the ℓ_∞ -norm and the ℓ_2 -norm of the pairwise-distortion vector using a grid search over the angular orientation of the line that represents the embedding. We evaluate the near-neighbor preservation of a given query point in the

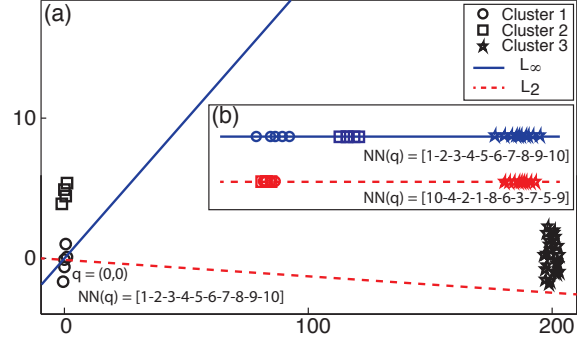


Figure 1. Comparison of the near-neighbor (NN) preservation performance of hashing based on minimizing the ℓ_∞ -norm (worst-case distortion) vs. the ℓ_2 -norm (average distortion) of the pairwise distance distortion vector on an illustrative data set. (a) For a dataset with three clusters (5 circles, 5 squares, and 60 stars), we found the optimal embeddings for both error metrics using grid search. (b) The projection of the data points using the ℓ_∞ -optimal embedding preserves three well-separated clusters; however the projection using the ℓ_2 -optimal embedding mixes circles and squares, projecting them into a single cluster. For the query point $q = (0, 0)$, all of its nearest neighbors $NN(q)$ are preserved with the correct ordering using the worst-based distortion embedding but not the average distortion embedding.

embedded space (shown in Fig. 1 (b)). For a query point q , located without loss of generality at the origin, the nearest neighbor ranking from the ambient space is destroyed by the ℓ_2 -optimal embedding, since the circle and square clusters overlap. In contrast, the ℓ_∞ -optimal embedding exactly preserves the rankings.

This illustration emphasizes the importance of preserving *relevant* distances in data retrieval tasks. To minimize the average distortion, the ℓ_2 -optimal embedding (dashed red line) sacrifices the square–circle distances in favor of the square–star and circle–star distances, which contribute more to the ℓ_2 -norm of the pairwise distance distortion vector. In contrast, the ℓ_∞ -optimal embedding focuses on the hardest distances to preserve (i.e., the worst-case distortion), leading to an embedding with smaller isometry constant than the ℓ_2 -optimal embedding. Preservation of these distances is critical for near-neighbor retrieval.

1.1. Contributions

We make four distinct contributions in this paper. First, conceptually, we advocate minimizing the worst-case distortion, which is formulated as an ℓ_∞ -norm minimization problem, and show that this approach outperforms approaches based on minimizing the average, ℓ_2 -norm distortion in a range of computer vision and learning scenarios (Hartley & Schaffalitzky, 2004).

Second, algorithmically, since ℓ_∞ -norm minimization problems are computationally challenging, especially for large datasets, we develop two accelerated and scalable al-

gorithms to find the optimal worst-case embedding. The first, *near-isometric binary hashing* (NIBH), is based on the alternating direction method of multipliers (ADMM) framework (Boyd et al., 2011). The second, NIBH-CG, is based on an accelerated greedy extension of the NIBH algorithm using the concept of *column generation* (Dantzig & Wolfe, 1960). NIBH-CG can rapidly learn hashing functions from large-scale data sets that require the preservation of *billions* of pairwise distances (e.g., *MNIST*).

Third, theoretically, since current data-dependent hashing algorithms do not offer any probabilistic guarantees in terms of preserving near-neighbors, we develop new theory to prove that, under natural assumptions regarding the data distribution and with a notion of hardness of near-neighbor search, NIBH preserves the nearest neighbors with high probability. Our analysis approach could be of independent interest for obtaining theoretical guarantees for other data-dependent hashing schemes.

Fourth, experimentally, we demonstrate the superior performance of NIBH as compared to ten state-of-the-art binary hashing algorithms using an exhaustive set of experimental evaluations involving six diverse datasets and three different performance metrics (near-isometry, Hamming distance ranking, and kendall τ ranking performance). In particular, we show that NIBH achieves the same distance preservation and Hamming ranking performance as state-of-the-art algorithms *while using up to 60% fewer bits*. Our experiments clearly show the superiority of the ℓ_∞ -norm formulation over the more classical ℓ_2 -norm formulation that underlies many hashing algorithms, such as BRE and IsoHash. Our formulation also outperforms recently developed techniques that assume more structure in their hash functions, such as *Spherical Hamming Distance Hashing* (SHD) (Heo et al., 2012) and *Circulant Binary Embedding* (CBE) (Yu et al., 2014).

2. Near-Isometric Binary Hashing

The standard formulation for data dependent binary hash function embeds a data point $\mathbf{x} \in \mathbb{R}^N$ into the low-dimensional Hamming space $\mathcal{H} = \{0, 1\}^M$ by first multiplying it by an *embedding matrix* $\mathbf{W} \in \mathbb{R}^{M \times N}$ and then quantizing the entries of the product $\mathbf{W}\mathbf{x}$ to binary values:

$$h(\mathbf{W}\mathbf{x}) = \frac{1 + \text{sgn}(\mathbf{W}\mathbf{x})}{2}. \quad (1)$$

The function $\text{sgn}(\cdot)$ operates element-wise on the entries of $\mathbf{W}\mathbf{x}$, transforming the real-valued vector $\mathbf{W}\mathbf{x}$ into a set of binary codes depending on the sign of the entries in $\mathbf{W}\mathbf{x}$.

2.1. Problem formulation

Consider the design of an embedding f that maps Q high-dimensional data vectors $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_Q\}$ in

the ambient space \mathbb{R}^N into low-dimensional binary codes $\mathcal{H} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_Q\}$ in the Hamming space with $\mathbf{h}_i \in \{0, 1\}^M$, where $\mathbf{h}_i = f(\mathbf{x}_i)$, $i = 1, \dots, Q$, and $M \ll N$. Define the distortion of the embedding by

$$\delta = \inf_{\lambda > 0} \sup_{(i,j) \in \Omega} |\lambda d_H(\mathbf{h}_i, \mathbf{h}_j) - d(\mathbf{x}_i, \mathbf{x}_j)|,$$

with $\Omega = \{(i, j) : i, j \in \{1, 2, \dots, Q\}, i > j\}$,

where $d(\mathbf{x}_i, \mathbf{x}_j)$ denotes the Euclidean distance between the data points $\mathbf{x}_i, \mathbf{x}_j$, $d_H(\mathbf{h}_i, \mathbf{h}_j)$ denotes the Hamming distance between the binary codes \mathbf{h}_i and \mathbf{h}_j , and λ is a positive scaling variable. The distortion δ measures the worst-case deviation from perfect isometry (i.e., optimal distance preservation) among all pairs of data points. Define the *secant set* $\mathcal{S}(\mathcal{X})$ as $\mathcal{S}(\mathcal{X}) = \{\mathbf{x}_i - \mathbf{x}_j : (i, j) \in \Omega\}$, i.e., the set of all pairwise difference vectors in \mathcal{X} . Let $|\mathcal{S}(\mathcal{X})| = |\Omega| = Q(Q-1)/2$ denote the size of the secant set. Note that the common distortion measure utilized in other hashing algorithms is the average distortion, i.e., $\sum_{i>j} (\lambda d_H(\mathbf{h}_i, \mathbf{h}_j) - d(\mathbf{x}_i, \mathbf{x}_j))^2 / |\Omega|$.

We formulate the problem of minimizing the distortion parameter δ as the following optimization problem:

$$\text{minimize}_{\mathbf{W}, \lambda > 0} \sup_{(i,j) \in \Omega} \left| \lambda \|h(\mathbf{W}\mathbf{x}_i) - h(\mathbf{W}\mathbf{x}_j)\|_2^2 - \|\mathbf{x}_i - \mathbf{x}_j\|_2 \right|,$$

since the squared ℓ_2 -distance between a pair of binary codes is equivalent to their Hamming distance up to a scaling factor that can be absorbed into λ . We can rewrite the above optimization problem as

$$(\mathbf{P}^*) \quad \text{minimize}_{\mathbf{W}, \lambda > 0} \|\lambda \mathbf{v}'(\mathbf{W}) - \mathbf{c}\|_\infty,$$

where $\mathbf{v}' \in \mathbb{R}^{Q(Q-1)/2}$ is a vector containing the pairwise Hamming distances between the embedded data vectors $\|h(\mathbf{x}_i) - h(\mathbf{x}_j)\|_2^2$, and \mathbf{c} is a vector containing the pairwise ℓ_2 -distances between the original data vectors. Intuitively, the ℓ_∞ -norm objective optimizes the *worst-case* distortion among all pairs of data points.

The problem (\mathbf{P}^*) is a combinatorial problem with complexity $\mathcal{O}(Q^{2M})$. To overcome the combinatorial nature of the problem, we approximate the hash function $h(\cdot)$ by the sigmoid function (also known as the inverse logit link function) $\sigma(x) = (1 + e^{-x})^{-1}$. This enables us to approximate (\mathbf{P}^*) by the following optimization problem:

$$(\mathbf{P}) \quad \text{minimize}_{\mathbf{W}, \lambda > 0} \|\lambda \mathbf{v}(\mathbf{W}) - \mathbf{c}\|_\infty,$$

where $\mathbf{v} \in \mathbb{R}_+^{Q(Q-1)/2}$ is a vector containing the pairwise ℓ_2 distances between the embedded data vectors after sigmoid relaxation $\|(1 + e^{-\mathbf{W}\mathbf{x}_i})^{-1} - (1 + e^{-\mathbf{W}\mathbf{x}_j})^{-1}\|_2^2$. Here the sigmoid function operates element-wise on $\mathbf{W}\mathbf{x}_i$. In practice we use a more general definition of the sigmoid

function, defined as $\sigma_\alpha(x) = (1 + e^{-\alpha x})^{-1}$, where α is the rate parameter controlling how closely it approximates the non-smooth function $h(\cdot)$. The following lemma characterizes the quality of such an approximation (see the Appendix for a proof).

Lemma 1. *Let x be a Gaussian random variable as $x \sim \mathcal{N}(\mu, \sigma^2)$. Define the distortion of the sigmoid approximation at x as $|h(x) - \sigma_\alpha(x)|$. Then, the expected distortion is bounded as $\mathbb{E}_x[|h(x) - \sigma_\alpha(x)|] \leq \frac{1}{\sigma\sqrt{2\pi\alpha}} + 2e^{-(\sqrt{\alpha}+c/\alpha\sigma^2)}$, where c is a positive constant. As α goes to infinity, the expected distortion goes to 0.*

Remark. As has been noted in the machine vision literature (Zoran & Weiss, 2012), a natural model for an image database is that its images are generated from a mixture of Gaussian distributions. Lem. 3 bounds the deviation of the sigmoid approximation from the non-smooth hash function (1) under this model.

2.2. Near-isometry and nearest neighbor preservation

Inspired by the definition of *relative contrast* in (He et al., 2012), we define a more generalized measure of data separability to preserve k -NN that we call the k -order gap $\Delta_k := d(\mathbf{x}_0, \mathbf{x}_{k+1}) - d(\mathbf{x}_0, \mathbf{x}_k)$, where \mathbf{x}_0 is a query point and \mathbf{x}_k and \mathbf{x}_{k+1} are its k^{th} and $k+1^{\text{th}}$ nearest neighbors, respectively. We formally show that if the data is highly separable (Δ_k is large), then the above approach preserves all k nearest neighbors with high probability (see the Appendix for a proof and discussion).

Theorem 2. *Assume that all the data points are independently generated from a mixture of Gaussian distribution i.e., $\mathbf{x}_i \sim \sum_{p=1}^P \pi_p \mathcal{N}(\mu_p, \Sigma_p)$. Let $\mathbf{x}_0 \in \mathbb{R}^N$ denote a query data point in the ambient space, and the other data points \mathbf{x}_i be ordered so that $d(\mathbf{x}_0, \mathbf{x}_1) < d(\mathbf{x}_0, \mathbf{x}_2) < \dots < d(\mathbf{x}_0, \mathbf{x}_Q)$. Let δ denote the final value of the distortion parameter computed from any binary hashing algorithm, and let c denote a positive constant. Then, if $\mathbb{E}_x[\Delta_k] \geq 2\delta + \sqrt{\frac{1}{c} \log \frac{Qk}{\epsilon}}$, the binary hashing algorithm preserves all the k -nearest neighbors of a data point with probability at least $1 - \epsilon$.*

2.3. The NIBH algorithm

We now develop an algorithm to solve the optimization problem (P). We apply the alternating direction method of multipliers (ADMM) framework (Boyd et al., 2011) to construct an efficient algorithm to find a (possibly local) optimal solution of (P). Note that (P) is non-convex, and therefore no standard optimization method is guaranteed to converge to a globally optimal solution in general. We introduce an auxiliary variable \mathbf{u} to arrive at the equivalent

problem:

$$\underset{\mathbf{W}, \mathbf{u}, \lambda > 0}{\text{minimize}} \quad \|\mathbf{u}\|_\infty \quad \text{subject to} \quad \mathbf{u} = \lambda \mathbf{v}(\mathbf{W}) - \mathbf{c}. \quad (2)$$

The augmented Lagrangian form of this problem can be written as $\underset{\mathbf{W}, \mathbf{u}, \lambda > 0}{\text{minimize}} \quad \|\mathbf{u}\|_\infty + \frac{\rho}{2} \|\mathbf{u} - \lambda \mathbf{v}(\mathbf{W}) + \mathbf{c} + \mathbf{y}\|_2^2$, where ρ is the scaling parameter in ADMM and $\mathbf{y} \in \mathbb{R}^{Q(Q-1)/2}$ is the Lagrange multiplier vector. The NIBH algorithm proceeds as follows. First, the variables \mathbf{W} , λ , \mathbf{u} , and Lagrange multipliers \mathbf{y} are initialized randomly. Then, at each iteration, we optimize over each of the variables \mathbf{W} , \mathbf{u} , and λ while holding the other variables fixed. More specifically, in iteration ℓ , we perform the following four steps until convergence:

- *Optimize over \mathbf{W} via $\mathbf{W}^{(\ell+1)} \leftarrow \arg \min_{\mathbf{W}} \frac{1}{2} \sum_{(i,j) \in \Omega} (u_{ij}^{(\ell)} - \lambda^{(\ell)} \|\frac{1}{1+e^{-\mathbf{W}\mathbf{x}_i}} - \frac{1}{1+e^{-\mathbf{W}\mathbf{x}_j}}\|_2^2 + \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 - y_{ij}^{(\ell)})^2$, where $\lambda^{(\ell)}$ denotes the value of λ in the ℓ^{th} iteration. We also use $u_{ij}^{(\ell)}$ and $y_{ij}^{(\ell)}$ to denote the entries in $\mathbf{u}^{(\ell)}$ and $\mathbf{y}^{(\ell)}$ that correspond to the pair \mathbf{x}_i and \mathbf{x}_j in the dataset \mathcal{X} . We show in our experiments below that using the accelerated first-order gradient descent algorithm (Nesterov, 2007) to solve this subproblem results in good empirical convergence performance (see the Appendix).*
- *Optimize over \mathbf{u} while holding the other variables fixed; it corresponds to solving the proximal problem of the ℓ_∞ -norm $\mathbf{u}^{(\ell+1)} \leftarrow \arg \min_{\mathbf{u}} \|\mathbf{u}\|_\infty + \frac{\rho}{2} \|\mathbf{u} - \lambda^{(\ell)} \mathbf{v}^{(\ell+1)} + \mathbf{c} + \mathbf{y}^{(\ell)}\|_2^2$. We use the low-cost algorithm described in (Studer et al., 2014) to perform the proximal operator update.*
- *Optimize over λ while holding the other variables fixed; it corresponds to a positive least squares problem, where λ is updated as $\lambda^{(\ell+1)} \leftarrow \arg \min_{\lambda > 0} \frac{1}{2} \|\mathbf{u}^{(\ell+1)} - \lambda \mathbf{v}^{(\ell+1)} + \mathbf{c} + \mathbf{y}^{(\ell)}\|_2^2$. We perform this update using the non-negative least squares algorithm (Kim & Park, 2007).*
- *Update \mathbf{y} via $\mathbf{y}^{(\ell+1)} \leftarrow \mathbf{y}^{(\ell)} + \eta(\mathbf{u}^{(\ell+1)} - \lambda^{(\ell+1)} \mathbf{v}^{(\ell+1)} + \mathbf{c})$, where the parameter η controls the dual update step size.*

2.4. Accelerated NIBH for large-scale datasets

The ADMM-based NIBH algorithm is efficient for small-scale datasets (e.g., for secant sets of size $|\mathcal{S}(\mathcal{X})| < 5000$ or so). However, the memory requirement of NIBH is quadratic in $|\mathcal{X}|$, which would be problematic for applications involving large-scale numbers of data points and secants. In response, we develop an algorithm that approximately solves (P) while scaling very well to large-scale problems. The key idea comes from classical results

in optimization theory related to *column generation* (CG) (Dantzig & Wolfe, 1960; Hegde et al., 2015).

The optimization problem (2) is an ℓ_∞ -norm minimization problem with an equality constraint on each secant. The Karush-Kuhn-Tucker (KKT) condition for this problem states that, if strong duality holds, then the optimal solution is entirely specified by a (typically very small) portion of the constraints. Intuitively, the secants corresponding to these constraints are the pairwise distances that are harder to preserve in the low-dimensional Hamming space. We call the set of such secants the *active set*. In order to solve (P), it suffices to find the active secants and solve NIBH with a much smaller number of active constraints. To leverage the idea of the active set, we iteratively run NIBH on a small subset of all the secants that violate the near-isometry condition, as detailed below:

- Solve (P) with a small random subset S_0 of all the secants $\mathcal{S}(\mathcal{X})$ using NIBH to obtain $\widehat{\mathbf{W}}$, $\hat{\delta}$, and $\hat{\lambda}$, initial estimates of the parameters. Identify the active set \mathcal{S}_a . Fix $\lambda = \hat{\lambda}$ for the rest of the algorithm.
- Randomly select a new subset $\mathcal{S}_v \subset \mathcal{S}$ of secants that violate the near isometry condition using the current estimates of $\widehat{\mathbf{W}}$, $\hat{\delta}$, and $\hat{\lambda}$. Then, form an augmented secant set $\mathcal{S}_{\text{aug}} = \mathcal{S}_a \cup \mathcal{S}_v$.
- Solve (P) with the secants in the set \mathcal{S}_{aug} using the NIBH algorithm.

We dub this approach *NIBH-CG*. NIBH-CG iterates over the above steps until no new violating secants are added to the active set. Since the algorithm searches over all the secants for violating secants in each iteration before terminating, NIBH-CG ensures that all of the constraints are satisfied when it terminates. A key benefit of NIBH-CG is that only the set of active secants (and not all secants) needs to be stored in memory. This benefit leads to significant improvements in terms of memory complexity over competing algorithms, since the set of all secants quickly becomes large-scale and can exceed the system memory capacity in large-scale applications.

3. Experiments

In this section, we validate the NIBH and NIBH-CG algorithms via experiments using a range of synthetic and real-world datasets, including three small-scale, three medium-scale, and one large-scale datasets with respect to three metrics. We compare NIBH against ten state-of-the-art binary hashing algorithms, including binary reconstructive embedding (BRE) (Kulis & Darrell, 2009), spectral hashing (SH) (Weiss et al., 2009), anchor graph hashing

(AGH) (Liu et al., 2011), multidimensional spectral hashing (MDSH) (Weiss et al., 2012), scalable graph hashing (SGH) (Jiang & Li, 2015), PCA hashing (PCAH) (Jolliffe, 2002), isotropic hashing (IsoHash) (Kong & Li, 2012), spherical Hamming distance hashing (SHD) (Heo et al., 2012), circulant binary embedding (CBE) (Yu et al., 2014), and locality-sensitive hashing (LSH) (Indyk & Motwani, 1998).

3.1. Performance metrics and datasets

We compare the algorithms using the following three metrics:

Maximum distortion $\delta = \inf_{\lambda > 0} \|\lambda \hat{\mathbf{v}} - \mathbf{c}\|_\infty$, where the vector $\hat{\mathbf{v}}$ contains the pairwise Hamming distances between the learned binary codes. This metric quantifies the distance preservation among all of the pairwise distances after projecting the training data in the ambient space into binary codes. We also define the maximum distortion for unseen test data δ_{test} , which measures the distance preservation on a hold-out test dataset using the hash function learned from the training dataset.

Mean average precision (MAP) for near-neighbor preservation in the Hamming space. MAP is computed by first finding the set of k -nearest neighbors for each query point on a hold-out test data in the ambient space \mathcal{L}^k and the corresponding set \mathcal{L}_H^k in the Hamming space and then calculating the average precision $\text{AP} = |\mathcal{L}^k \cap \mathcal{L}_H^k|/k$. We then report MAP by calculating the mean value of AP across all data points.

Kendall τ ranking correlation coefficient. We first rank the set of k -nearest neighbors for each data point by increasing distance in the ambient space as $\mathcal{T}(\mathcal{L}^k)$ and in the Hamming space as $\mathcal{T}(\mathcal{L}_H^k)$. The Kendall τ correlation coefficient is a scalar $\tau \in [-1, 1]$ that measures the similarity between the two ranked sets $\mathcal{T}(\mathcal{L}^k)$ and $\mathcal{T}(\mathcal{L}_H^k)$ (Kendall, 1938). The value of τ increases as the similarity between the two rankings increases and reaches the maximum value of $\tau = 1$ when they are identical. We report the average value of τ across all data points in the training dataset.

To compare the algorithms, we use the following standard datasets from computer vision: *Random* consists of independently drawn random vectors in \mathbb{R}^{100} from a multivariate Gaussian distribution with zero mean and identity covariance matrix. *Translating squares* is a synthetic dataset consisting of 10×10 images that are translations of a 3×3 white square on black background (Hegde et al., 2015). *MNIST* is a collection of 60,000 28×28 grayscale images of handwritten digits (LeCun & Cortes, 1998). *Photo-Tourism* is a corpus of approximately 300,000 image patches, represented using scale-invariant feature transform (SIFT) features (Lowe, 2004) in \mathbb{R}^{128} (Snavely et al.,

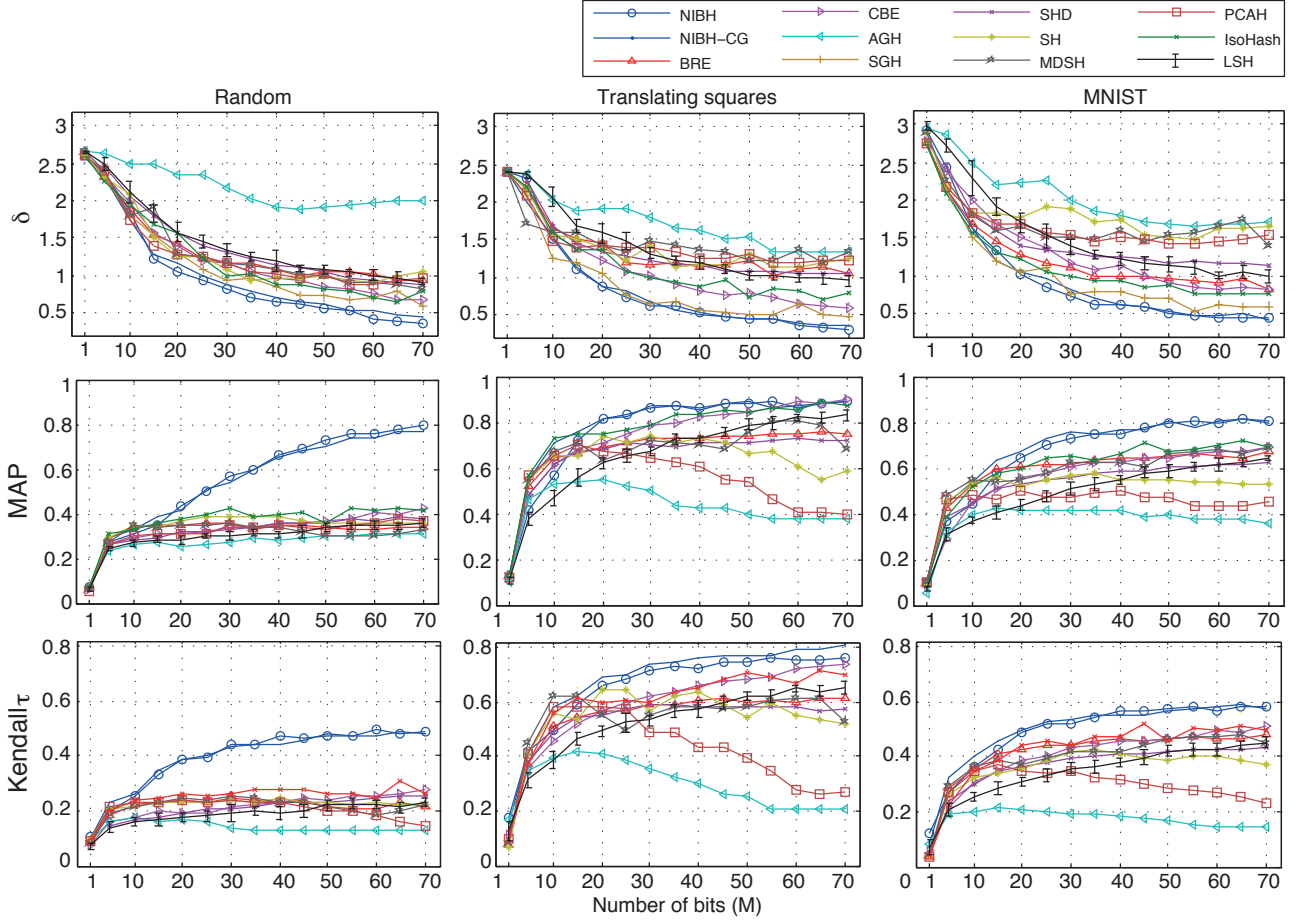


Figure 2. Comparison of the NIBH and NIBH-CG algorithms against several baseline binary hashing algorithms using three small-scale datasets with 4950 secants ($Q = 100$). The performance of NIBH-CG closely follows that of NIBH, and both outperform all of the other algorithms in terms of the maximum distortion δ (superior distance preservation), mean average precision MAP of training samples (superior nearest neighbor preservation), and Kendall τ rank correlation coefficient (superior ranking preservation).

2006). *LabelMe* is a collection of over 20,000 images represented using GIST descriptors in \mathbb{R}^{512} (Torralba et al., 2008). *Peekaboom* is a collection of 60,000 images represented using GIST descriptors in \mathbb{R}^{512} (Torralba et al., 2008). Following the experimental approaches of the hashing literature (Kulis & Darrell, 2009; Norouzi & Fleet, 2011), we pre-process the data by subtracting the mean and then normalizing all points to lie on the unit sphere.

3.2. Small- and medium-scale experiments

We start by evaluating the performance of NIBH and NIBH-CG using a small-scale subset of the first three datasets. Small-scale datasets enable us to compare the performance of NIBH vs. NIBH-CG to verify that they perform similarly. Also they help us assess the asymptotic behavior of algorithms in preserving isometry since the total of number of secants are small compare to the bit budget in compact binary codes.

Experimental setup. We randomly select $Q = 100$ data points from the *Random*, *Translating squares*, and *MNIST* datasets. We then apply the NIBH, NIBH-CG, and all the baseline algorithms on each dataset for different target binary code word lengths M from 1 to 70 bits. We set the NIBH and NIBH-CG algorithm parameters to the common choice of $\rho = 1$ and $\eta = 1.6$. To generate hash function of length M for LSH, we draw M random vectors from a Gaussian distribution with zero mean and an identity covariance matrix. We use the same random vectors to initialize NIBH and other baseline algorithms. In the near-neighbor preservation experiments, to show the direct advantage of minimizing ℓ_∞ -norm over ℓ_2 -norm, we followed the exact procedure described in BRE (Kulis & Darrell, 2009) to select the training secants, i.e., we apply the NIBH algorithm on only the lowest 5% of the pairwise distances (which are set to zero as in BRE) combined with the highest 2% of the pairwise distances.

We follow the *continuation* approach (Wen et al., 2010) to

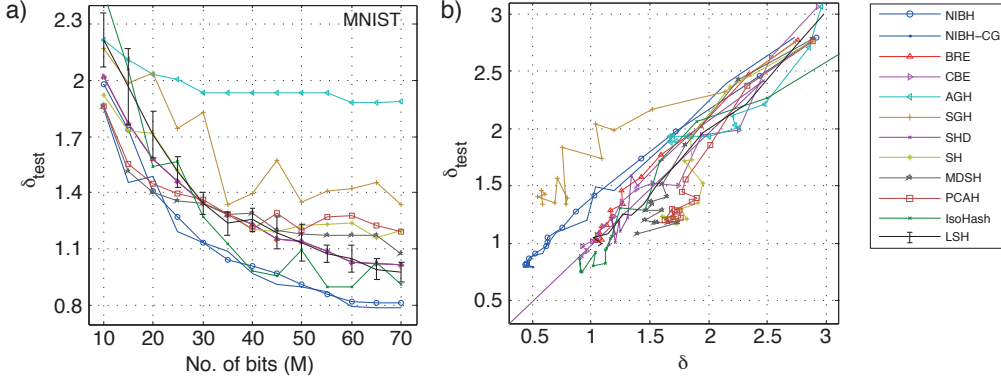


Figure 3. Comparison of NIBH and NIBH-CG against several state-of-the-art binary hashing algorithms in preserving isometry on MNIST data. (a) NIBH-CG outperforms the other algorithms in minimizing the isometry constant on unseen data δ_{test} . (b) NIBH and NIBH-CG provide better isometry guarantee with a small sacrifice to universality.

set the value of α . We start with a small value of α , (e.g., $\alpha = 1$) to avoid becoming stuck in bad local minima, and then gradually increase α as the algorithm proceeds. As the algorithm gets closer to convergence and has obtained a reasonably good estimate of the parameters \mathbf{W} and λ , we set $\alpha = 10$, which enforces a good approximation of the sign function (see Lemma 3).

Results. The plots in the top row of Figure 2 illustrate the value of the distortion parameter δ as a function of the number of projections (bits) M . The performance of NIBH and NIBH-CG closely follow each other, indicating that NIBH-CG is a good approximation to NIBH. Both NIBH and NIBH-CG outperform the other baseline algorithms in terms of the distortion parameter δ . Among these baselines, LSH has the lowest isometry performance since random projections are oblivious to the intrinsic geometry of the training dataset. To achieve $\delta = 1$, NIBH(-CG) requires 60% fewer bits M than CBE and BRE. NIBH(-CG) also achieves better isometry performance asymptotically, i.e., up to $\delta \approx 0.5$, given a sufficient number of bits ($M \geq 70$), while for most of the other algorithms the performance plateaus after $\delta = 1$. NIBH’s superior near-isometry performance extends well to unseen data. Figure 3(a) demonstrates that NIBH achieves the lowest isometry constant on a test dataset δ_{test} compared to other hashing algorithms. Figure 3(b) further suggests that NIBH’s superior isometry performance comes with smallest sacrifice to the universality of the hash functions.

The plots in the middle and bottom row of Figure 2 shows the average precision for retrieving training data and the Kendall τ correlation coefficient respectively, as a function of the number of bits M . We see that NIBH preserves a higher percentage of nearest neighbors compared to other baseline algorithms as M increases with better average ranking among $k = 10$ -nearest neighbors.

Now we showcase the performance of NIBH-CG on three medium-scale, real-world datasets used in (Kulis & Darrell, 2009; Norouzi & Fleet, 2011), including *Photo-tourism*, *LabelMe*, and *Peekaboom* for the popular machine learning task of *data retrieval*. From each dataset we randomly select $Q = 1000$ training points, following the setup in BRE (Kulis & Darrell, 2009), and use them to train NIBH-CG and the other baseline algorithms. We then randomly select a separate set of $Q = 1000$ data points and use it to test the performance of NIBH-CG and other baseline algorithms in terms MAP with $k = 50$. Figure 4 illustrates the performance of NIBH-CG on these datasets. NIBH-CG outperforms all the baseline algorithms with large margins in Hamming ranking performance in term of MAP with top-50 near-neighbors.

3.3. Large-scale experiments

We now demonstrate that NIBH-CG scales well to large-scale datasets. We use the full *MNIST* dataset with 60,000 training images and augment it with three rotated versions of each image (rotations of 90° , 180° , and 270°) to create a larger dataset with $Q = 240,000$ data points. Next, we construct 4 training sets with 1,000, 10,000, 100,000, and 240,000 images out of this large set. We train all algorithms with $M = 30$ bits and compare their performance on a test set of 10,000 images. BRE fails to execute on a standard desktop PC with 12 GB of RAM for training sets with more than 100,000 points due to the size of the secant set $|\mathcal{X}|$. The results for all algorithms are given in Table 1; we tabulate their performance in terms of MAP for the top-500 neighbors. The performance of NIBH-CG is significantly better than the baseline algorithms and, moreover, improves as the size of the training set grows. This emphasizes that NIBH-CG excels at large-scaled problems thanks to its very small memory requirement; indeed, the memory requirement of NIBH-CG is linear in the number of *active*

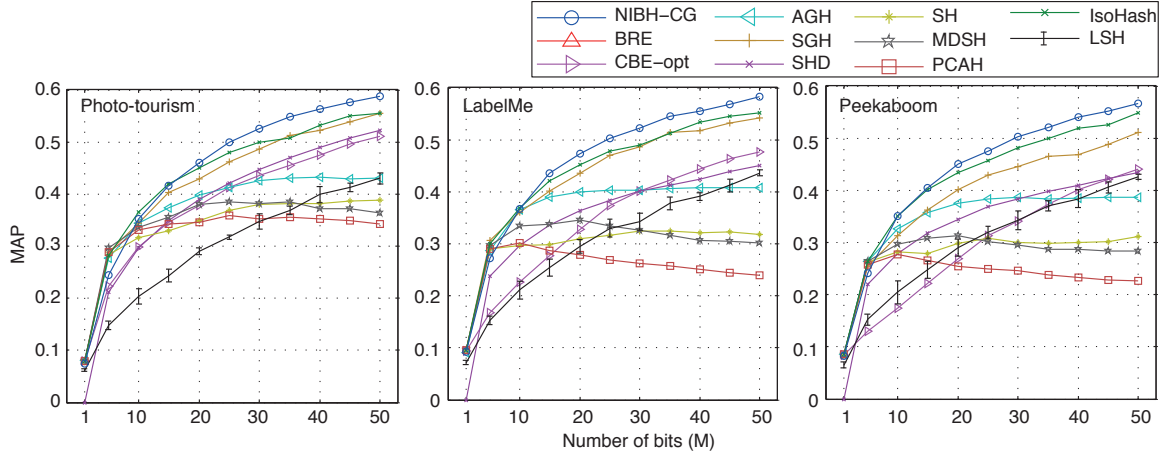


Figure 4. Hamming ranking performance comparison on three medium-scale datasets ($Q = 1000$). The top-50 neighbors are used to report MAP over a test data of same size.

secants rather than the total number of secants.

4. Discussion

We have demonstrated that the worst-case, ℓ_∞ -norm-based near-isometric binary hashing (NIBH) algorithm is superior to a wide range of algorithms based on the more traditional average-case, ℓ_2 -norm. Despite its non-convexity and non-smoothness, NIBH admits an efficient optimization algorithm that converges to a high-performing local minimum. Moreover, NIBH-CG, the accelerated version of NIBH, provides significant memory advantages over existing algorithms. Our exhaustive experiments with six datasets, three metrics, and ten algorithms have shown that NIBH outperforms all of the state-of-the-art data-dependent hashing algorithms. The results in this paper provide a strong motivation for exploring ℓ_∞ -norm formulations in binary hashing.

References

- Andoni, A., Indyk, P., Nguyen, H. L., and Razenshteyn, I. Beyond locality-sensitive hashing. In *Proc. 25th Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 1018–1028, Jan 2014.
- Boyd, S., Parikh, N., Chu, E., Peleato, B., and Eckstein, J. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, Jan. 2011.
- Cover, T. M. and Thomas, J. A. *Elements of information theory*. John Wiley & Sons, 2012.
- Dantzig, G. B. and Wolfe, P. Decomposition principle for linear programs. *Operations Research*, 8(1):101–111, Feb. 1960.
- Datar, M., Immorlica, N., Indyk, P., and Mirrokni, V. S. Locality-sensitive hashing scheme based on p-stable distributions. In *Proc. 20th Annual Symposium on Computational Geometry*, pp. 253–262, June 2004.
- Gong, Y. and Lazebnik, S. Iterative quantization: A procrustean approach to learning binary codes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 817–824, June 2011.
- Hartley, R. and Schaffalitzky, F. L_∞ minimization in geometric reconstruction problems. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 504–509, June 2004.
- He, J., Kumar, S., and Chang, S. On the difficulty of nearest neighbor search. In *Proc. 29th Intl. Conf. on Machine Learning*, pp. 1127–1134, June 2012.
- Hegde, C., Sankaranarayanan, A. C., Yin, W., and Baraniuk, R. G. Numax: A convex approach for learning near-isometric linear embeddings. *IEEE Trans. Signal Processing*, 63(22):6109–6121, Nov. 2015.
- Heo, J., Lee, Y., He, J., Chang, S., and Yoon, S. Spherical hashing. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 2957–2964, 2012.
- Horn, R. A. and Johnson, C. R. *Topics in Matrix Analysis*. Cambridge University Press, 1991.
- Indyk, P. and Motwani, R. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proc. 30th Annual ACM Symposium on Theory of Computing*, pp. 604–613, May 1998.
- Jiang, Q. Y. and Li, W. J. Scalable graph hashing with feature transformation. *IJCAI*, 2015.
- Jolliffe, I. *Principal Component Analysis*. Wiley Online Library, 2002.

Table 1. Comparison of NIBH-CG against several baseline binary hashing algorithms on large-scale *MNIST* datasets with over 28 billion secants $|\mathcal{S}(\mathcal{X})|$. We tabulate the Hamming ranking performance in terms of mean average precision MAP for different sizes of the dataset. All training times are in seconds.

$M = 30$ bits	MAP / Top-500 (<i>MNIST</i> + <i>rotations</i>)				Training time
Training size Q	1K	10K	100K	240K	240K
Secant size $ \mathcal{S}(\mathcal{X}) $	500K	50M	5B	28B	28B
NIBH-CG	52.79 (± 0.15)	54.69 (± 0.18)	54.93 (± 0.23)	55.52 (± 0.11)	541.43
BRE	48.33 (± 0.65)	50.67 (± 0.33)	—	—	18685.51
CBE	38.70 (± 1.18)	38.12 (± 1.34)	38.50 (± 2.05)	38.53 (± 0.83)	68.94
SPH	44.33 (± 0.74)	44.24 (± 0.61)	44.37 (± 0.71)	44.32 (± 0.63)	184.46
SH	40.12 (± 0.00)	39.37 (± 0.00)	38.79 (± 0.00)	38.26 (± 0.00)	3.05
MDSH	41.06 (± 0.00)	41.23 (± 0.00)	40.80 (± 0.00)	40.39 (± 0.00)	15.00
AGH	45.81 (± 0.34)	47.78 (± 0.38)	47.69 (± 0.41)	47.38 (± 0.32)	4.49
SGH	51.32 (± 0.07)	51.33 (± 0.20)	51.01 (± 0.23)	50.66 (± 0.76)	5.89
PCAH	39.90 (± 0.00)	38.53 (± 0.00)	38.81 (± 0.00)	37.50 (± 0.00)	0.08
IsoHash	50.91 (± 0.00)	50.90 (± 0.00)	50.72 (± 0.00)	50.55 (± 0.00)	2.82
LSH	33.69 (± 0.94)	33.69 (± 0.94)	33.69 (± 0.94)	33.69 (± 0.94)	2.29×10^{-4}

- Kendall, M. G. A new measure of rank correlation. *Biometrika*, 30(1–2):81–93, June 1938.
- Kim, H. and Park, H. Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis. *Bioinformatics*, 23 (12):1495–1502, Apr. 2007.
- Kong, W. and Li, W. J. Isotropic hashing. In *Advances in Neural Information Processing Systems*, pp. 1646–1654, 2012.
- Kulis, B. and Darrell, T. Learning to hash with binary reconstructive embeddings. In *Advances in Neural Information Processing Systems*, pp. 1042–1050, Dec. 2009.
- LeCun, Y. and Cortes, C. The MNIST database of handwritten digits, 1998.
- Ledoux, M. and Talagrand, M. *Probability in Banach Spaces: Isoperimetry and Processes*. Springer-Verlag, 1991.
- Liu, W., Wang, J., Kumar, S., and Chang, S. Hashing with graphs. In *Proc. 28th intl. Conf. on Machine Learning*, pp. 1–8, July 2011.
- Lowe, D. G. Distinctive image features from scale-invariant keypoints. *Intl. Journal of Computer Vision*, 60(2):91–110, Nov. 2004.
- Lv, Q., Josephson, W., Wang, Z., Charikar, M., and Li, K. Multi-probe lsh: efficient indexing for high-dimensional similarity search. In *Proceedings of the 33rd international conference on Very large data bases*, pp. 950–961. VLDB Endowment, 2007.
- Nesterov, Y. Gradient methods for minimizing composite objective function. Technical report, Université Catholique de Louvain, Sep. 2007.
- Norouzi, M. and Fleet, D. J. Minimal loss hashing for compact binary codes. In *Proc. 28th Intl. Conf. on Machine Learning*, pp. 353–360, June 2011.
- Plan, Y. and Vershynin, R. Dimension reduction by random hyperplane tessellations. *Discrete & Computational Geometry*, 51(2):438–461, Mar. 2014.
- Shaw, B. and Jebara, T. Minimum volume embedding. In *Proc. 11th Intl. Conf. on Artificial Intelligence and Statistics*, pp. 460–467, Mar. 2007.
- Snaveley, N., Seitz, S. M., and Szeliski, R. Photo tourism: Exploring photo collections in 3D. *ACM Trans. Graphics*, 25(3):835–846, July 2006.
- Studer, C., Goldstein, T., Yin, W., and Baraniuk, R. G. Democratic representations. *Preprint*, 2014. URL <http://www.csl.cornell.edu/~studer/papers/14TIT-linf-submitted.pdf>.
- Torralba, A., Fergus, R., and Weiss, Y. Small codes and large image databases for recognition. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 1–8, June 2008.
- Wainwright, M. J. Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (Lasso). *IEEE Trans. Information Theory*, 55(5):2183–2202, May 2009.
- Wang, J., Shen, H. T., Song, J., and Ji, J. Hashing for similarity search: A survey. *arXiv preprint arXiv:1408.2927*, 2014.
- Weinberger, K. Q. and Saul, L. K. Unsupervised learning of image manifolds by semidefinite programming. *Intl. Journal of Computer Vision*, 70(1):77–90, May 2006.

- Weiss, Y., Torralba, A., and Fergus, R. Spectral hashing. In *Advances in Neural Information Processing Systems*, pp. 1753–1760, Dec. 2009.
- Weiss, Y., Fergus, R., and Torralba, A. Multidimensional spectral hashing. In *European Conf. on Computer Vision*, pp. 340–353. Oct. 2012.
- Wen, Z., Yin, W., Goldfarb, D., and Zhang, Y. A fast algorithm for sparse reconstruction based on shrinkage, subspace optimization, and continuation. *SIAM Journal on Scientific Computing*, 32(4):1832–1857, June 2010.
- Yu, F. X., Kumar, S., Gong, Y., and Chang, S. Circulant binary embedding. In *Proc. 31st Intl. Conf. on Machine Learning*, pp. 946–954, June 2014.
- Zoran, D. and Weiss, Y. Natural images, gaussian mixtures and dead leaves. In *Advances in Neural Information Processing Systems*, pp. 1736–1744, Dec. 2012.

A. Appendix

In the appendix, we prove Lem. 3 and Thm. 4 on the performance of NIBH on k -nearest neighbor preservation. Finally, we include additional numerical simulation results and discussions on the empirical convergence of the NIBH algorithm.

Proof of Lem. 3

Lemma 3. *Let x be a Gaussian random variable as $x \sim \mathcal{N}(\mu, \sigma^2)$. Define the distortion of the sigmoid approximation at x as $|h(x) - \sigma_\alpha(x)|$. Then, the expected distortion is bounded as*

$$\mathbb{E}_x[|h(x) - \sigma_\alpha(x)|] \leq \frac{1}{\sigma\sqrt{2\pi\alpha}} + 2e^{-(\sqrt{\alpha}+c/\alpha\sigma^2)},$$

where c is a positive constant. As α goes to infinity, the expected distortion goes to 0.

Proof. It is easy to see that the distortion $|h(x) - \sigma_\alpha(x)|$ occurs at $x = 0$. Therefore, among different values of μ , $\mu = 0$ gives the largest distortion since the density of x peaks at $x = 0$. Therefore, we bound the distortion at setting $\mu = 0$, which is an upper bound of the distortion when $\mu \neq 0$. By definition (1) in the main text, $h(x)$ can be written as

$$h(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

When $x \sim \mathcal{N}(0, \sigma^2)$, we have

$$\begin{aligned} \mathbb{E}_x[|h(x) - \sigma_\alpha(x)|] &= \int_{-\infty}^{\infty} |h(x) - \sigma_\alpha(x)| \mathcal{N}(x; 0, \sigma^2) dx \\ &= 2 \int_0^{\infty} (h(x) - \sigma_\alpha(x)) \mathcal{N}(x; 0, \sigma^2) dx \\ &= 2 \int_0^{x_0} (h(x) - \sigma_\alpha(x)) \mathcal{N}(x; 0, \sigma^2) dx \\ &\quad + 2 \int_{x_0}^{\infty} (h(x) - \sigma_\alpha(x)) \mathcal{N}(x; 0, \sigma^2) dx \\ &\leq 2 \int_0^{x_0} \frac{1}{2} \mathcal{N}(x; 0, \sigma^2) dx + 2 \int_{x_0}^{\infty} \frac{1}{1 + e^{\alpha x_0}} \mathcal{N}(x; 0, \sigma^2) dx \end{aligned} \quad (3)$$

$$\leq \frac{x_0}{\sqrt{2\pi}\sigma} + 2 \frac{e^{-cx_0^2/\sigma^2}}{1 + e^{\alpha x_0}} \quad (5)$$

$$\leq \frac{1}{\sigma\sqrt{2\pi\alpha}} + 2e^{-(\sqrt{\alpha}+c/\alpha\sigma^2)}, \quad (6)$$

when we set $x_0 = \frac{1}{\sqrt{\alpha}}$ and c is a positive constant. In (3), we used the fact that $\sigma_\alpha(x)$ and $h(x)$ are symmetric with respect to the point $(0, \frac{1}{2})$. (4) is given by the properties of the sigmoid function, (5) is given by the Gaussian concentration inequality (Ledoux & Talagrand, 1991), and (6) is

given by the inequality $1/(1 + e^{\alpha x_0}) \leq e^{-\alpha x_0}$. The fact that $\mathbb{E}_x[|h(x) - \sigma_\alpha(x)|] \rightarrow 0$ as $\alpha \rightarrow \infty$ is obvious from the bound above. \square

Proof of Thm. 4

Theorem 4. Assume that all the data points are independently generated from a mixture of Gaussian distribution, i.e., $\mathbf{x}_i \sim \sum_{p=1}^P \pi_p \mathcal{N}(\mu_p, \Sigma_p)$. Let $\mathbf{x}_0 \in \mathbb{R}^N$ denote a query data point in the ambient space, and the other data points \mathbf{x}_i be ordered so that $d(\mathbf{x}_0, \mathbf{x}_1) < d(\mathbf{x}_0, \mathbf{x}_2) < \dots < d(\mathbf{x}_0, \mathbf{x}_Q)$. Let δ denote the final value of the distortion parameter computed from any binary hashing algorithm, and let c denote a positive constant. Then, if $\mathbb{E}_x[\Delta_k] \geq 2\delta + \sqrt{\frac{1}{c} \log \frac{Qk}{\epsilon}}$, the binary hashing algorithm preserves the k -nearest neighbors of a point with probability at least $1 - \epsilon$.

In order to prove this theorem, we need the following Lemma:

Lemma 5. Let $\mathbf{x}_0, \dots, \mathbf{x}_N$ and Δ_k be defined as in Thm. 4. Then, there exist a constant c such that $P(\Delta_k - \mathbb{E}_x[\Delta_k] < t) \leq e^{-ct^2}$ for $t > 0$.

Proof. Since the data points \mathbf{x}_0 , \mathbf{x}_k and \mathbf{x}_{k+1} are independently generated from a finite mixture of Gaussian distributions, the random variable of their concatenation $\mathbf{y} = [\mathbf{x}_0^T, \mathbf{x}_k^T, \mathbf{x}_{k+1}^T]^T \in \mathbb{R}^{3N}$ is sub-Gaussian (Wainwright, 2009). Then, we have

$$\begin{aligned} \Delta_k(\mathbf{y}) &= \|\mathbf{x}_0 - \mathbf{x}_{k+1}\|_2 - \|\mathbf{x}_0 - \mathbf{x}_k\|_2 \\ &= \left\| \begin{pmatrix} \mathbf{I} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -\mathbf{I} \end{pmatrix} \mathbf{y} \right\|_2 - \left\| \begin{pmatrix} \mathbf{I} & 0 & 0 \\ 0 & -\mathbf{I} & 0 \\ 0 & 0 & 0 \end{pmatrix} \mathbf{y} \right\|_2 \\ &\leq \left\| \begin{pmatrix} 2\mathbf{I} & 0 & 0 \\ 0 & -\mathbf{I} & 0 \\ 0 & 0 & -\mathbf{I} \end{pmatrix} \mathbf{y} \right\|_2 \\ &\leq 2\|\mathbf{y}\|_2, \end{aligned}$$

where we have used the triangular inequality in the second to last step, the Rayleigh-Ritz theorem (Horn & Johnson, 1991), and the fact that the maximum singular value of the matrix in the step before is 2. This result means that $\Delta_k(\mathbf{y})$ is a Lipschitz function of \mathbf{y} . Thus, by Talagrand's inequality (Ledoux & Talagrand, 1991), we have that $P(\Delta_k - \mathbb{E}_x[\Delta_k] < -t) \leq e^{-ct^2}$ for some positive constant c and $t > 0$, since \mathbf{y} is sub-Gaussian. \square

Now we are ready to prove Thm. 4.

Proof. Let E denote the event that the set of top- k nearest neighbors is not preserved in the Hamming space. Then, we have $E = \cup e_{m,n}$, where $e_{m,n}$ denote the event that

$d_H(\mathbf{x}_0, \mathbf{x}_m) > d_H(\mathbf{x}_0, \mathbf{x}_n)$ with $m \in \{1, \dots, k\}$ and $n \in \{k+1, \dots, Q\}$. Then, using the union bound (Cover & Thomas, 2012), we have

$$\begin{aligned} P(E) &\leq \sum_{m,n} P(e_{m,n}) \leq k(Q-k)P(e_{k,k+1}) \\ &= k(Q-k)P(d_H(\mathbf{x}_0, \mathbf{x}_k) > d_H(\mathbf{x}_0, \mathbf{x}_{k+1})) \\ &= k(Q-k)P(d_H(\mathbf{x}_0, \mathbf{x}_{k+1}) < d_H(\mathbf{x}_0, \mathbf{x}_k)), \end{aligned}$$

where we have used the fact that the most possible event among all $e_{m,n}$ events is the one corresponding to the order mismatch between the k^{th} and $k+1^{\text{th}}$ nearest neighbor. Now, note that the NIBH output δ satisfies $\max_{i,j} |d_H(\mathbf{x}_i, \mathbf{x}_j) - d(\mathbf{x}_i, \mathbf{x}_j)| \leq \delta^1$. Observe that $\Delta_k = d(\mathbf{x}_0, \mathbf{x}_{k+1}) - d(\mathbf{x}_0, \mathbf{x}_k) \geq 2\delta$ is a sufficient condition for $d_H(\mathbf{x}_0, \mathbf{x}_{k+1}) \geq d_H(\mathbf{x}_0, \mathbf{x}_k)$, since

$$\begin{aligned} d_H(\mathbf{x}_0, \mathbf{x}_{k+1}) - d_H(\mathbf{x}_0, \mathbf{x}_k) &\geq d(\mathbf{x}_0, \mathbf{x}_{k+1}) - \delta - d(\mathbf{x}_0, \mathbf{x}_k) - \delta \\ &\geq 2\delta - 2\delta = 0, \end{aligned}$$

by the triangular inequality. This leads to

$$\begin{aligned} P(d_H(\mathbf{x}_0, \mathbf{x}_{k+1}) < d_H(\mathbf{x}_0, \mathbf{x}_k)) &= 1 - P(d_H(\mathbf{x}_0, \mathbf{x}_{k+1}) \geq d_H(\mathbf{x}_0, \mathbf{x}_k)) \\ &\leq 1 - P(\Delta_k \geq 2\delta) = P(\Delta_k < 2\delta). \end{aligned}$$

Therefore, combining all the above and Lem. 5, the probability that the k -nearest neighbor is not preserved is bounded by

$$\begin{aligned} P(E) &\leq k(Q-k)P(d_H(\mathbf{x}_0, \mathbf{x}_{k+1}) < d_H(\mathbf{x}_0, \mathbf{x}_k)) \\ &\leq k(Q-k)P(\Delta_k < 2\delta) \\ &= k(Q-k)P(\Delta_k - \mathbb{E}_x[\Delta_k] < -(\mathbb{E}_x[\Delta_k] - 2\delta)) \\ &\leq k(Q-k)e^{-c(\mathbb{E}_x[\Delta_k] - 2\delta)^2} \\ &\leq kQe^{-c(\mathbb{E}_x[\Delta_k] - 2\delta)^2}. \end{aligned}$$

Now, let $kQe^{-c(\mathbb{E}_x[\Delta_k] - 2\delta)^2} \leq \epsilon$, we have that the requirement for the k -nearest neighbors to be exactly preserved with probability at least $1 - \epsilon$ is

$$\mathbb{E}_x[\Delta_k] \geq 2\delta + \sqrt{\frac{1}{c} \log \frac{Qk}{\epsilon}}.$$

\square

Remark Note that our bound on the number of nearest neighbor preserved k depends on the final outcome of the NIBH algorithm in the value of δ . In order to relate our result to the number of binary hash functions M required, we can make use of (Plan & Vershynin, 2014,

¹Here we assume $\lambda = 1$ without loss of generality.

Thm. 1.10). For a bounded set $K \subset \mathbb{R}^N$ with diameter 1, let $M \geq C\delta^{-6}w(K)^2$ where $w(K) = \mathbb{E}_x[\sup_{\mathbf{x} \in K} \langle g, \mathbf{x} \rangle]$ denotes the Gaussian width of K and some constant C . Then, (Plan & Vershynin, 2014, Thm. 1.10) states that, with high probability, $h(x)$ as defined in (1) with a random matrix \mathbf{W} whose entries are independently generated from $\mathcal{N}(0, 1)$ is a δ_0 -isometric embedding. Therefore, if we initialize NIBH with such a random \mathbf{W} which is likely to be δ_0 -isometric, then empirically (see Figure 5), the NIBH algorithm will learn a better embedding that is δ -isometric with $\delta < \delta_0$. Therefore, we have that the number of hash functions M required for k -nearest neighbor preservation is at least $M \sim (\mathbb{E}_x[\Delta_k] - \sqrt{\log(kQ/\epsilon)}^{-6}w(\mathbf{X})^2$ assuming that the training dataset \mathbf{X} is properly normalized.

Empirical convergence of the NIBH algorithm

Figure 5 shows the empirical loss and the actual distortion parameter δ as a function of the iteration count ℓ in the NIBH algorithm as applied on 4950 secants (i.e., $Q = 100$) from the *MNIST* dataset. The behavior of the empirical loss function closely matches that of δ as they gradually converge. The curve empirically confirms that minimizing the loss function in each iteration of NIBH (using the ADMM framework) directly penalizes the non-convex loss function (distortion parameter δ). After initializing the NIBH algorithm with random entries for \mathbf{W} , the value of the distortion parameter δ significantly drops in the first few iterations, empirically confirming that NIBH learns an embedding with significantly lower distortion parameter after a few iterations.

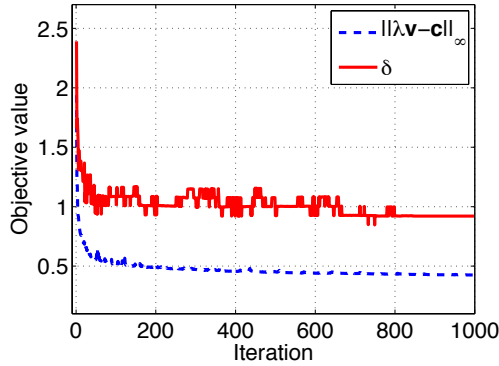


Figure 5. Empirical convergence behavior of the NIBH algorithm. Both the maximum distortion parameter δ and the loss function $\|\lambda \mathbf{v} - \mathbf{c}\|_\infty$ that approximates δ gradually decrease and converge as the number of iterations increases. We see that the loss function of NIBH closely matches the behavior of the actual distortion parameter in each iteration.